# An automatic robust Bayesian approach to principal component regression

Philippe Gagnon [1], Mylène Bédard [2], Alain Desgagné [3]

January 28, 2020

[1]Department of Statistics, University of Oxford, United Kingdom.
[2]Department of Mathematics and Statistics, Université de Montréal, Canada.
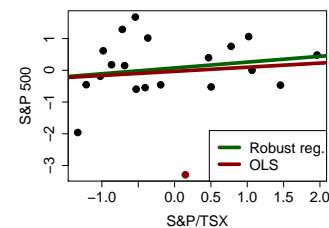[3]Department of Mathematics, Université du Québec à Montréal, Canada.

### Abstract

Principal component regression uses principal components as regressors. It is particularly useful in prediction settings with high-dimensional covariates. The existing literature treating of Bayesian approaches is relatively sparse. We introduce a Bayesian approach that is robust to outliers in both the dependent variable and the covariates. Outliers can be thought of as observations that are not in line with the general trend. The proposed approach automatically penalises these observations so that their impact on the posterior gradually vanishes as they move further and further away from the general trend, corresponding to a concept in Bayesian statistics called *whole robustness*. The predictions produced are thus consistent with the bulk of the data. The approach also exploits the geometry of principal components to efficiently identify those that are significant. Individual predictions obtained from the resulting models are consolidated according to model-averaging mechanisms to account for model uncertainty. The approach is evaluated on real data and compared to its nonrobust Bayesian counterpart, the traditional frequentist approach, and a commonly employed robust frequentist method. Detailed guidelines to automate the entire statistical procedure are provided. All required code is made available, see ArXiv:1711.06341.

Keywords: dimension reduction; linear regression; outliers; principal component analysis; reversible jump algorithms; whole robustness.

## 1 Introduction

In statistical analyses, information carried by several variables is commonly summarised to allow visualisation or model estimation when the number of variables makes it unstable or impossible. For instance, S&P 500 and S&P/TSX respectively summarise the stock prices of 500 and about 250 large companies domiciled in the United States and Canada, and are commonly used to portray the American and Canadian economies. Figure 1 illustrates the relationship between the January 2011 daily returns of these two indices. The scatter plot is further summarised using two different linear regression



**Figure 1.** January 2011 daily returns

models that respectively yield a robust regression line (in green) and an ordinary least squares regression line (in red). Given that different summaries (other than S&P 500 or S&P/TSX in our example) lead to different data points and therefore different regressions, one might however wonder whether the available or natural summaries are necessarily suitable for the tasks at hand.

Principal component regression (PCR) is the name given to a linear regression model using principal components (PCs) as regressors. It is based on a principal component analysis (PCA), which is commonly used to summarise the information contained in covariates. The principle is to find new axes in the covariate space by exploiting the correlation structure between the covariates, and then encode the covariate observations in that new coordinate system. The resulting variables, called principal components (PCs), are linearly independent and have the remarkable property that the first $q$ PCs retain the maximum amount of information carried by the original observations (compared to any other $q$-dimensional summary). Regrouping correlated variables to produce linearly independent ones is appealing in a linear regression context, as strongly correlated variables are known to carry redundant information, leading to unstable estimates. Companies within the same economic sector in stock market indices like S&P 500 and S&P/TSX are an example of such correlated variables. Linear independence also allows visualising the relationship between the dependent variable and the PCs by plotting the dependent variable against each of the PCs.

Due to the loss in the interpretability of the inference results engendered by transforming covariates, PCR is mainly used in a prediction context. It can nevertheless be useful for clarifying the underlying structure in the original covariates, as shown in West (2003). In this paper, we consider a Bayesian prediction framework and address four issues; those are described below.

**Robustness against outliers.** It is common knowledge that OLS (ordinary least squares) estimates become significantly contaminated in presence of outliers. In Figure 1, the OLS regression line (in red) is pulled below the robust regression line (in green) by the outlier (red dot). OLS estimates make the assumption that errors are normally distributed, which affects both the linear regression and the PCA (see Šmídl and Quinn (2007)). In presence of outliers, the slimness of normal tails causes a shift in the posterior so as to incorporate the information carried by all the data. The posterior may thus find itself concentrated in an area that is not supported by any source of information, be it the outliers or the bulk of the data. This translates, for instance, into predictions that are not in line with either of these two groups.

The natural solution to this problem is to assume an error distribution with heavier tails, and therefore more adapted to the possible presence of outliers. The Student distribution becomes an obvious choice as it leads to a straightforward implementation of the Bayesian regression approach via the Gibbs sampler (West, 1984). Using a heavy-tailed distribution like the Student however only allows attaining partial robustness (Andrade and O'Hagan, 2011), which may lead to regression coefficients with inflated variances, and ultimately contaminated model selection. Relying on an uncontaminated model selection procedure is crucial in our framework as the identification of important PCs relies on it.

It was recently proved in Gagnon *et al.* (2018) that model selection in linear regression is uncontaminated when a super heavy-tailed error distribution is instead assumed. We follow this path, and based on that strategy of using super heavy-tailed distributions, introduce a new class of wholly robust Bayesian PCA. Hereafter, whole robustness refers to an approach that automatically penalises observations that are not in line with the general trend, so that their impact on the posterior distribution gradually vanishes as they move further from that trend. The assumed super heavy-tailed density matches the standard normal outside of the tails, which makes the approaches efficient. The resemblance between the two densities helps us to design the computational tools.

**Selection of significant PCs.** The selection of pertinent PCs to be included in our robust regression model is based on model selection and in line with the methods used in Wang (2012) and Tipton *et al.* (2017). Ours however differs in that we do not use the stochastic search variable selection (see George and McCulloch (1993)), which is the common tool to discriminate among a large number of (typically correlated) regressors. We instead take advantage of the linear independence among PCs to quickly exclude the irrelevant ones, leading to the following two-step approach. We first evaluate the individual relevance of each PC through Bayes factors, after which the retained PCs are used to propose a sequence of nested models. The joint posterior of these models and their parameters is next computed. Observations for the dependent variable are predicted by accounting for model uncertainty through model averaging (see, for instance, Raftery *et al.* (1997) and Hoeting *et al.* (1999)).

**Automatic and efficient implementation.** Our approach to attain whole robustness (which consists in assuming super heavy-tailed error distributions) however prevents us from having access to full conditional distributions and, therefore, to using Gibbs sampler. For the robust PCA, we then propose a simplified computational scheme based on point estimates. The model posterior probabilities are however required in the linear regression stage of the statistical analysis, and so we turn to the reversible jump (RJ) algorithm to obtain estimates of these probabilities. The RJ sampler is a Markov chain Monte Carlo (MCMC) method introduced by Green (1995) that allows to directly sample from the joint posterior of the models and their parameters. The efficiency of such samplers relies heavily on the design of the functions required for the implementation. We provide a detailed procedure to automatically implement an efficient RJ algorithm.

**Prior specification.** It is often difficult, in PCR, to specify meaningful priors on the models and their parameters. For this reason, noninformative priors are commonly favoured. The simplest noninformative structure is arguably the improper Jeffreys priors on the parameters of all models, along with a uniform prior on the models. With such a prior structure, one might wonder whether the so-called Jeffreys-Lindley paradox (Lindley, 1957; Jeffreys, 1967), representing inconsistent model selection results, may arise. We show that this is not the case and adopt that structure.

**Structure of the paper.** The general model is described in Section 2. Nonrobust normal PCA and regression approaches are presented in Section 3, followed by their robust counterparts, representing the proposed methodology, in Section 4. In particular, the proposed robust PCA is discussed in Section 4.1, while the robust linear regression is addressed in Section 4.2. Section 4.2.1 presents the RJ sampler and then Section 4.2.2 focuses on automating its implementation. The stock market indices example is revisited in Section 5 where all the features of the proposed robust approach are illustrated. The validity of our prior structure is addressed in the supplementary material (Section 7) as this part is not required to understand and implement the proposed methodology.

## 2 Principal component regression

Consider that we have access to a rank $r \in \{1, 2, \ldots\}$ matrix $\mathbf{C} \in \mathbb{R}^{n \times p}$ containing $n \in \{1, 2, \ldots\}$ observations from $p \in \{1, 2, \ldots\}$ standardised covariates. A PCA is then performed on this data set. It will be seen that standardisation and PCA in the proposed robust approach are different from those in its nonrobust counterpart. We thus defer details about these steps to later sections.

Denote by $\mathbf{Z}_q$ the matrix of rank $q \leq r$ arising from either dimension reduction technique (nonrobust or robust PCA). The design matrix $\mathbf{X} := (x_{ij})$ is constructed by simply grafting a column vector of 1's to the matrix $\mathbf{Z}_q$. For simplicity, we will refer to this extra column of $\mathbf{X}$ as the first component. The PCs are thus contained in the following columns, and $d := q + 1$ denotes the number of columns of $\mathbf{X}$.

We wish to study the relationship between a dependent variable with data points $Y_1, \ldots, Y_n \in \mathbb{R}$ and the PCs in order to predict values for the former. We start from the premise that the relationship is linear:

$$Y_i = \mathbf{x}_{i,K}^T \boldsymbol{\beta}_K + \epsilon_{i,K}, \quad i = 1, \ldots, n, \ \ K \in \{1, \ldots, \mathrm{K}_{\max}\}, \tag{1}$$

where $K$ is the model indicator, $\mathrm{K}_{\max}$ is a positive integer representing the number of models considered, and $\epsilon_{1,K}, \ldots, \epsilon_{n,K} \in \mathbb{R}$ are the errors associated to Model $K$. The vector of observed PCs included in Model $K$ satisfies $\mathbf{x}_{i,K} := \{x_{ij} : j \in I_K\}$, where $I_k \subseteq \{1, \ldots, d\}$ is a vector whose elements indicate which PCs are included in Model $K = k$. For instance, $I_1$ is associated to Model 1 which, in this paper, always corresponds to the model containing only the intercept ($I_1 := \{1\}$). The $d_K$-dimensional vector of regression coefficients associated to Model $K$ is $\boldsymbol{\beta}_K := (\beta_{1,K}, \ldots, \beta_{d_K,K})^T \in \mathbb{R}^{d_K}$, where $d_K$ is the cardinality of $I_K$. As is typically done in Bayesian linear regression, we assume that $\epsilon_{1,K}, \ldots, \epsilon_{n,K}$ and $\boldsymbol{\beta}_K$ are $n + 1$ conditionally independent random variables given $(K, \sigma_K)$, with $\sigma_K > 0$ being the scale parameter of the errors of Model $K$. The conditional density of $\epsilon_{i,K}$ is given by

$$\epsilon_{i,K} \mid K, \sigma_K, \boldsymbol{\beta}_K \overset{\mathcal{D}}{=} \epsilon_{i,K} \mid K, \sigma_K \overset{\mathcal{D}}{\sim} (1/\sigma_K) f(\epsilon_{i,K}/\sigma_K), \quad i = 1, \ldots, n.$$

Even though we assume a linear relationship between the dependent variable and regressors in (1), we remain realistic and adopt George Box's point of view, which says that all models are wrong, but that some are useful. The degree of usefulness represented by the model fits will presumably be reflected in the posterior model probabilities.

To study the relationship between the dependent variable and the PCs, we first identify the statistically relevant PCs. The individual contribution of the various components is assessed using Bayes factors. Specifically, we consider in the first step of the statistical analysis the $d$ models associated to $I_1 := \{1\}, I_2 := \{1, 2\}, \ldots, I_d = \{1, d\}$, and compare each of Models 2 through $d$ to Model 1. The PCs associated to Bayes factors greater than a given threshold are retained in the second step of the statistical analysis; the others are discarded.

In the second step of the analysis, we consider the sequence of nested models arising from the statistically significant PCs and find the posterior probabilities of these models, along with their parameter estimates. For instance, if the first, second and fourth PCs are the only ones deemed relevant, the sequence of models is $I_1 := \{1\}, I_2 := \{1, 2\}$, and $I_3 := \{1, 2, 4\}$. Considering only a sequence of nested models is natural in our context, as PCA generates components that carry less and less information about the original covariates; that also simplifies subsequent computations.

Finding posterior probabilities and parameter estimates is achieved by sampling from the joint posterior distribution of $(K, \sigma_K, \boldsymbol{\beta}_K)$ given $\mathbf{y} := (y_1, \ldots, y_n)^T$, denoted by $\pi(k, \sigma_k, \boldsymbol{\beta}_k \mid \mathbf{y})$, where the domain of $k$ depends on which step of the analysis is performed (and, for the second step, on the results of the previous step). Once estimates are obtained in the second step, values for the dependent variable can be predicted through model-averaging mechanisms.

# 3 Normal nonrobust models

## 3.1 Traditional principal component analysis

Several strategies allow retrieving the usual PCA from estimates of statistical models (see, e.g., Tipping and Bishop (1999) and Šmídl and Quinn (2007)). These methods assume that $\mathbf{C}$ has been generated from a linear model with normal errors. One can thus view PCs as point estimates and conduct a full Bayesian analysis of the model. We follow here the approach of Šmídl and Quinn (2007); its

presentation will facilitate the introduction of the robust PCA model as it will be analogously defined in Section 4.1.

The singular value decomposition allows expression of the matrix $\mathbf{C}$ as $\mathbf{ZLA}^T$, where $\mathbf{A}$ is a $p \times r$ matrix whose columns are the eigenvectors $\mathbf{v}_1, \ldots, \mathbf{v}_r$ of the sample correlation matrix of $\mathbf{C}$ with corresponding eigenvalues $\lambda_1 \geq \lambda_2 \geq \ldots \geq \lambda_r$, $\mathbf{L}$ is a $r \times r$ diagonal matrix with diagonal entries given by (up to a constant) $\lambda_1, \lambda_2, \ldots, \lambda_r$, and $\mathbf{Z}$ is a $n \times r$ matrix whose $j$-th column is given by $\lambda_j^{-1/2}\mathbf{Cv}_j$; see Jolliffe (2011) for instance. The PCs are traditionally defined as the vectors $\mathbf{Cv}_j$. We consider hereafter that the eigenvalues $\lambda_j$ are the sample variances of the PCs. The vectors $\lambda_j^{-1/2}\mathbf{Cv}_j$ therefore correspond to standardised PCs. Recall that the PCs are additionally pairwise orthogonal.

With $q < r$, let $\mathbf{Z}_q$ and $\mathbf{A}_q$ be the matrices comprised of the first $q$ columns of $\mathbf{Z}$ and $\mathbf{A}$, respectively, and $\mathbf{L}_q$ be the diagonal matrix with diagonal entries given by $\lambda_1, \ldots, \lambda_q$. If we want to further reduce the dimension of $\mathbf{Z}$ to $n \times q$, and therefore approximately reconstruct $\mathbf{C}$, Šmídl and Quinn (2007) present a model and a set of assumptions under which the maximum likelihood solution that arises is the anticipated matrix $\mathbf{Z}_q$. The model is

$$\mathbf{C} = \mathbf{M} + \mathbf{E}, \tag{2}$$

where $\mathbf{M}$ is assumed to have rank $q$ (and can therefore be decomposed using the singular value decomposition as above), and entries of $\mathbf{E}$ are assumed to be independently distributed as $\mathcal{N}(0, \eta^2)$, $\eta > 0$. The maximum likelihood estimate (MLE) of $\mathbf{M}$ is $\mathbf{Z}_q\mathbf{L}_q\mathbf{A}_q^T$. This follows from the fact that $\mathbf{Z}_q\mathbf{L}_q\mathbf{A}_q^T$ minimises the total squared reconstruction error among rank $q$ matrices. The MLE corresponds to the maximum a posteriori (MAP) estimate when the prior is flat. We use the matrix $\mathbf{Z}_q$ to form our design matrix $\mathbf{X}$ in the nonrobust linear regressions.

It usually is good practice to cap the percentage of the total variation that is accounted for as above a certain threshold, eigenvectors are essentially numerical noise. In the numerical analyses we limit it to 95%, meaning that $q$ is the maximum value such that $\sum_{j=1}^{q} \lambda_j / \sum_{j=1}^{r} \lambda_j \leq 0.95$.

**Remark 1.** *It is clear from (2) that $\mathbf{C}$ is viewed as a matrix containing observations from random variables. This may be confusing at first given that regressors are usually treated as known constants. In our case, the regressors are a function of $\mathbf{C}$; they are thus initially treated as observations from random variables in the PCA part of the statistical analysis. We next consider $\mathbf{Z}_q$ (or its robust version) as known constants in the regression part of the analysis. Our approach can thus be viewed as an approximation to the full and exact Bayesian analysis, in which all random unknown quantities, including $\mathbf{Z}_q, \boldsymbol{\beta}_K$, and $\sigma_K$, would be in linear models and estimated simultaneously, conditionally on $\mathbf{C}$ and $\mathbf{y}$. Our approach aims at simplifying the computation and interpretation of the statistical procedure.*

## 3.2   Ordinary least squares regression

Under the normality of the error distribution in the linear regressions (i.e. assuming that $f := \mathcal{N}(0, 1)$), the joint posterior $\pi(k, \sigma_k, \boldsymbol{\beta}_k \mid \mathbf{y})$ leads to closed-form expressions for the posterior model probabilities and parameter estimates. These expressions, detailed in Proposition 1 below, are handy for comparing the results arising from our robust approach to those obtained under the normality assumption in the numerical analyses. They will also be used in the design of the RJ algorithm to sample from the posterior under the super heavy-tailed distribution assumption. Indeed, the super heavy-tailed distribution that we use is similar to the normal distribution, except in the tails. When there is no outlier, this thus leads to a posterior that is similar to that under normality. In the presence of outliers, the full posterior of the robust model is similar to the posterior based on the nonoutliers only (i.e. excluding the outliers)

under normality. In either case, relying on the structure of the posterior under normality is therefore suitable for designing the RJ algorithm.

**Proposition 1.** *Assume that $f := \mathcal{N}(0, 1)$ and let the conditional prior density of $(\sigma_K, \boldsymbol{\beta}_K)$ given K be $\pi(\sigma_k, \boldsymbol{\beta}_k | k) \propto 1/\sigma_k$. Then, the posterior can be factorised as*

$$\pi(k, \sigma_k, \boldsymbol{\beta}_k \mid \mathbf{y}) = \pi(k \mid \mathbf{y}) \, \pi(\sigma_k \mid k, \mathbf{y}) \prod_{j=1}^{d_k} \pi(\beta_{j,k} \mid k, \sigma_k, \mathbf{y}),$$

*where $k \in \{1, \ldots, K_{\max}\}, \sigma_k > 0, \boldsymbol{\beta}_k \in \mathbb{R}^{d_k}$,*

$$\pi(k \mid \mathbf{y}) \propto \frac{\pi(k) \, \Gamma((n - d_k)/2) \, \pi^{d_k/2}}{\left( \|\mathbf{y} - \widehat{\mathbf{y}}_k\|_2^2 / (n - 1) \right)^{\frac{n - d_k}{2}}}, \tag{3}$$

$$\pi(\sigma_k \mid k, \mathbf{y}) = \frac{2^{1 - \frac{n - d_k}{2}} \left( \|\mathbf{y} - \widehat{\mathbf{y}}_k\|_2^2 \right)^{\frac{n - d_k}{2}}}{\Gamma((n - d_k)/2) \, \sigma_k^{n - d_k + 1}} \exp \left\{ -\frac{1}{2\sigma_k^2} \|\mathbf{y} - \widehat{\mathbf{y}}_k\|_2^2 \right\},$$

*$\beta_{1,K} \mid K, \sigma_K, \mathbf{y} \sim \mathcal{N}(\widehat{\beta}_{1,K} := 0, \sigma_K^2/n)$, and finally $\beta_{j,K} \mid K, \sigma_K, \mathbf{y} \sim \mathcal{N}(\widehat{\beta}_{j,K} := \sum_{i=1}^n x_{iI_{j,K}} y_i/(n-1), \sigma_K^2/(n-1))$ for $j = 2, \ldots, d_K$ (if $K \geq 2$). Here, $\| \cdot \|_2$ is the Euclidean norm, $\widehat{\mathbf{y}}_k := \mathbf{x}_{i,k}^T \widehat{\boldsymbol{\beta}}_k, \widehat{\boldsymbol{\beta}}_k := (\widehat{\beta}_{1,k}, \ldots, \widehat{\beta}_{d_k,k})^T$, $I_{j,K}$ is the j-th component of $I_K$, and $\pi(k)$ is the prior of K. Note that the normalisation constant of $\pi(k \mid \mathbf{y})$ is the sum over k of the expression on the right-hand side of (3).*

*Proof.* See the supplementary material (Section 7).     □

In our analyses, we use Bayesian model averaging to predict values for the dependent variable given sets of observations from the covariates. When normality is assumed, we can therefore use $\mathbb{E}[Y_{n+1} \mid \mathbf{y}] = \sum_k \pi(k \mid \mathbf{y}) \mathbf{x}_{n+1,k}^T \widehat{\boldsymbol{\beta}}_k$, where $\widehat{\boldsymbol{\beta}}_k$ is defined in Proposition 1. Note that under normality, $\sigma_K^2 \mid K, \mathbf{y}$ has an inverse-gamma distribution with shape and rate parameters given by $(n - d_K)/2$ and $\|\mathbf{y} - \widehat{\mathbf{y}}_K\|_2^2/2$, respectively.

# 4   Proposed robust models

The proposed solution to limit the impact of outliers in PCA and linear regression is simple: replace the traditional normality assumption on the error terms by a super heavy-tailed distribution assumption. The super heavy-tailed distribution that we use is the log-Pareto-tailed standard normal (LPTN) distribution with parameter $\rho \in (2\Phi(1) - 1, 1) \approx (0.6827, 1)$, where $\Phi$ is the cumulative distribution function of a standard normal. This distribution has been introduced in Desgagné (2015) and is expressed as

$$f(x) := \begin{cases} \varphi(x) & \text{if} \quad |x| \leq \tau, \\ \varphi(\tau) \frac{\tau}{|x|} \left( \frac{\log \tau}{\log |x|} \right)^{\lambda + 1} & \text{if} \quad |x| > \tau, \end{cases} \tag{4}$$

where $x \in \mathbb{R}$. The terms $\tau > 1$ and $\lambda > 0$ are functions of $\rho$ and satisfy

$$\tau := \Phi^{-1}((1 + \rho)/2) := \{\tau : \mathbb{P}(-\tau \leq Z \leq \tau) = \rho \text{ for } Z \overset{\mathcal{D}}{\sim} \mathcal{N}(0, 1)\},$$

$$\lambda := 2(1 - \rho)^{-1}\varphi(\tau)\,\tau\log(\tau),$$

with $\varphi(\cdot)$ and $\Phi^{-1}(\cdot)$ respectively being the probability density function (PDF) and inverse cumulative distribution function of a standard normal. The parameter $\rho$ controls the size of the interval over which $f$ exactly matches the standard normal density (i.e. the interval $[-\tau, \tau]$). Outside of this area, the tails behave according to a log-Pareto density $(1/|x|)(\log|x|)^{-\lambda-1}$, hence its name.

Setting $\rho$ to 0.95 has proved to be suitable for practical purposes, as addressed in Desgagné (2015) for location-scale models and in Gagnon *et al.* (2018) for linear regression. Accordingly, this is the value that will be used in our numerical analyses. Smaller values lead to improved robustness, but also to models that are further from normality (which then lead to discrepancies among estimations in the absence of outliers).

The theoretical result that motivates the use of super heavy-tailed distributions has been introduced in Gagnon *et al.* (2018). It establishes that, as outliers (because of extreme dependent and/or covariate observations) move further and further away from the general trend, the posterior distribution of $(K, \sigma_K, \beta_K)$ arising from the whole data set converges towards the posterior of $(K, \sigma_K, \beta_K)$ arising from the nonoutliers only. To prove this, it is however necessary to assume that there are at most $\lfloor n/2 - (\max d_k - 1/2) \rfloor$ outliers in the data set, with $\lfloor \cdot \rfloor$ being the floor function. For a fixed max $d_k$, this condition translates into a limiting breakdown point of 50% as $n \longrightarrow \infty$.

As explained in Gagnon *et al.* (2018), these models have built-in robustness that resolves conflict in a sensitive way. It takes full consideration of nonoutliers and excludes observations that are undoubtedly outlying; in between these two extremes, it balances and bounds the impact of possible outliers. In other words, there is no need to explicitly identify outliers; the method automatically deals with the level of (un)certainty about the nature of the observations (nonoutliers, clear outliers or potential outliers), which is particularly valuable in high-dimensional and model selection problems. The robust models and their properties are the subject of a whole article. For brevity purposes, we refer the interested reader to Gagnon *et al.* (2018) for more details.

## 4.1 Robust principal component analysis

Attempts at robustifying the traditional PCA model in (2) have been made by various authors (see, for instance, Luttinen *et al.* (2009) and Zhao *et al.* (2014)). They however follow the model specification of Tipping and Bishop (1999) as opposed to that of Šmídl and Quinn (2007) (as we do here), and accordingly do not explicitly impose a rank constraint on the matrix $\mathbf{M}$ used to reconstruct $\mathbf{C}$. As mentioned in Section 3.1, this constraint ensures that $\mathbf{M}$ can be decomposed as $\tilde{\mathbf{Z}}_q\tilde{\mathbf{L}}_q\tilde{\mathbf{A}}_q$, where $\tilde{\mathbf{Z}}_q$ and $\tilde{\mathbf{A}}_q$ have orthogonal columns (and are estimated by $\mathbf{Z}_q$ and $\mathbf{A}_q$ under the normal errors assumption). This orthogonality combined with the properties of PCA lead to the appealing geometric interpretation that those new axes are the best to reflect the information contained in $\mathbf{C}$. It also facilitates the statistical procedure for identifying relevant regressors. The price to pay for these advantages under the robust model is a significant increase in terms of computational complexity, as it becomes necessary to perform sampling and optimisation within the manifold of orthogonal matrices. As an alternative to this computationally demanding route, we propose here an asymptotic approximation to a wholly robust PCA (as $n \longrightarrow \infty$ and outliers move further away from the general trend). An exhaustive analysis of the exact version (including its implementation) will be conducted separately.

In wholly robust PCA, the entries of the error matrix $\mathbf{E} := (e_{ij})$ are such that $e_{ij} \mid \eta \overset{\mathcal{D}}{\sim} (1/\eta)g(e_{ij}/\eta)$, with $g$ the density of the LPTN. Under this error distribution assumption, we conjecture that a convergence result similar to that proved in Gagnon *et al.* (2018) holds. In particular, the posterior distribution

of $(\tilde{\mathbf{Z}}_q, \tilde{\mathbf{L}}_q, \tilde{\mathbf{A}}_q, \eta)$ (obtained from the covariate matrix $\mathbf{C}$ under LPTN errors) converges towards the posterior of $(\tilde{\mathbf{Z}}_q, \tilde{\mathbf{L}}_q, \tilde{\mathbf{A}}_q, \eta)$ obtained from a new covariate matrix $\mathbf{C}^*$ and LPTN errors, as the outliers move away from the trend. Generally speaking, $\mathbf{C}^*$ is a matrix in which outlying covariate observations are vertically projected onto a regression plane that is obtained using the nonoutliers only. The proposed approximation to a wholly robust PCA makes use of the fact that the model with LPTN errors is similar to that with normal errors for the same reasons as Section 3.2, and thus essentially consists in computing the PCs using $\mathbf{C}^*\mathbf{v}_j^*$ as in Section 3.1, with $\mathbf{v}_j^*$ being the $j$-th eigenvector of a robust correlation matrix of $\mathbf{C}^*$. To better understand what happens, we consider an example containing a single PC which is simple enough for the wholly robust PCA model to be estimated. The orthogonality is indeed trivially verified given that there is only one column in $\tilde{\mathbf{Z}}_q$ and $\tilde{\mathbf{A}}_q$.

Suppose that $\mathbf{C}$ is a $21 \times 2$ matrix of observed covariates. Observations from the first covariate are $c_{i1} = i - 11$, $i = 1, \ldots, 21$, and observed values from the second one are generated from the model $c_{i2} = c_{i1} + \epsilon_i$ with $\epsilon_i \sim \mathcal{N}(0, 1)$, $i = 1, \ldots, n$. Figure 2 (a) illustrates the relationship between the observed covariates.

Let us now introduce an outlier in this sample by moving $(c_{21,1}, c_{21,2}) = (10, 10.92)$ to $(10, 20)$; this sample is represented by the black dots in Figure 2 (b). Applying a traditional PCA to these observed covariates and then using it to retrieve the matrix $\mathbf{C}$ yield the red dots in Figure 2 (b); the reconstruction using the traditional PCA can be seen to rotate around the centre of the data as the outlier moves away from the trend. The wholly robust PCA approach leads to different results. The reconstruction of $\mathbf{C}$ using that approach is represented by the yellow dots in Figure 2 (b).
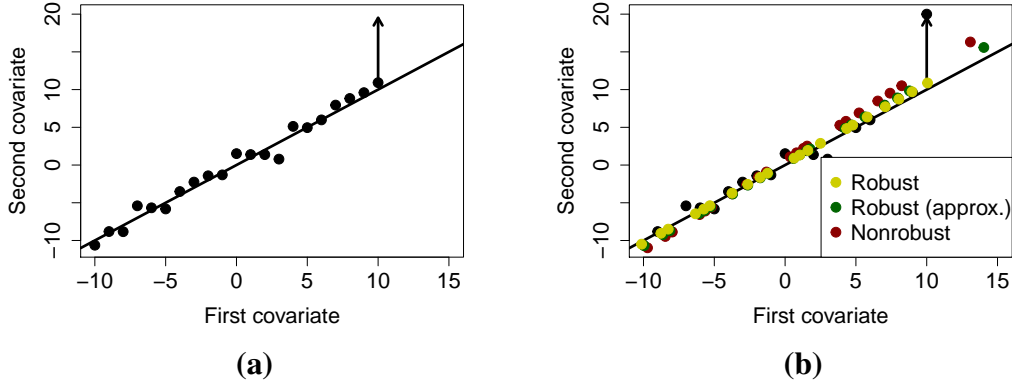
Now, suppose that the outlier $(c_{21,1}, c_{21,2})$ is vertically projected onto a regression line that is obtained using the first 20 observed covariates (i.e. the nonoutlying points only); in other words, the outlier is replaced by its predicted value at $c_{21,1}$. Denote this new covariate matrix by $\mathbf{C}^*$. It is observed that as $j$ increases in $(c_{21,2}, c_{21,2}) = (10, 10.92 + j)$, the posterior distribution of $(\tilde{\mathbf{Z}}_q, \tilde{\mathbf{L}}_q, \tilde{\mathbf{A}}_q, \eta)$ (obtained from $\mathbf{C}$, which includes the outlier) converges towards the posterior of $(\tilde{\mathbf{Z}}_q, \tilde{\mathbf{L}}_q, \tilde{\mathbf{A}}_q, \eta)$ obtained from $\mathbf{C}^*$. Applying the approximate robust PCA and then using it to reconstruct $\mathbf{C}$ yield the green dots in Figure 2 (b).

In that figure, it is seen that the exact and approximate robust approaches (yellow and green dots) produce very similar results; the reconstruction of the outlier is however different under both approaches (we explain why it is the case and why it is not a problem in robust PCR in the following paragraphs). It turns out that as the outlier $(c_{21,2}, c_{21,2}) = (10, 10.92 + j)$ reaches $(10, 20)$, the posterior distribution of $(\tilde{\mathbf{Z}}_q, \tilde{\mathbf{L}}_q, \tilde{\mathbf{A}}_q, \eta)$ (based on $\mathbf{C}$) has essentially converged. Indeed, moving the outlier further upwards has no effect on the results from the robust approaches; this is obviously not the case for the traditional PCA, which pursues its rotation around the centre of the data. For the data set in Figure 2 (b), the squared reconstruction errors based on the nonoutliers only are 8.77 and 16.39 for the approximate robust and nonrobust PCA, respectively; the exact robust method yields a similar result to its approximate counterpart.

We now detail the implementation of the approximate robust PCA.

1. Standardise the columns of the original data set to obtain $\mathbf{C}$ using the robust location-scale model of Desgagné (2015), with an LPTN error distribution and $\rho := 0.95$. This model corresponds to the linear regression model with the intercept only and $f := \text{LPTN}$. Location and scale estimates $\widehat{\mu}_j$ and $\widehat{\sigma}_j$ are thus used to standardise Column $j$, $j = 1, \ldots, p$.

2. Compute robust correlations between all pairs of columns in $\mathbf{C}$ using the slope estimator of the robust simple regression model with an LPTN error distribution and $\rho := 0.95$. These correlations form the robust correlation matrix. For simplicity, we set the upper diagonal entries to $\widehat{\beta}_{j_1, j_2}$,

**Figure 2.** (a) $n = 21$ points generated from the model $c_{i2} = c_{i1} + \epsilon_i$ with $\epsilon_i \sim \mathcal{N}(0, 1)$; (b) data set with outlier, and reconstruction using one PC under the robust and nonrobust PCA; the lines $y = x$ are also depicted

which denote the robust correlations between the standardised Columns $j_1$ and $j_2$ where Column $j_2 > j_1$ plays the role of the dependent variable. We next make the matrix symmetrical and set its diagonal elements to 1.

3. Compute the PCs $\widehat{\mathbf{Z}}_q$ using $\mathbf{C}\widehat{\mathbf{v}}_j$, with $\widehat{\mathbf{v}}_j$ being the $j$-th eigenvector of the robust correlation matrix of $\mathbf{C}$.

A mathematical justification of this approximation is presented in the supplementary material (Section 7). It is shown that $\mathbf{C}$ is asymptotically equivalent to $\mathbf{C}^*$ except for the components where there are outliers. Also, $\widehat{\mathbf{v}}_j$ is asymptotically equivalent to $\mathbf{v}_j^*$. Therefore there might be extreme values in the robust PCs, as there might be some in $\mathbf{C}$. If they exist, these extreme values will be handled by the robust linear regressions given their ability to deal with all types of outliers including leverage points.

Under the exact robust PCA approach, $\tilde{\mathbf{Z}}_q$ is directly estimated from the robust model; that represents the difference with the approximate method. The main advantage in using the approximate robust PCA is computational: the required estimates $\widehat{\mu}_j, \widehat{\sigma}_j$, and $\widehat{\beta}_{j_1,j_2}$ are easily obtained and can be computed in parallel. In our numerical experiments $\widehat{\mu}_j, \widehat{\sigma}_j$, and $\widehat{\beta}_{j_1,j_2}$ are maximum a posteriori (MAP) estimates with flat priors (corresponding to MLE). A second advantage is that the method allows automatic outlier detection. As in Gagnon *et al.* (2018), we compute estimates of the standardised residuals in the simple linear regressions as $z_i^{j_1,j_2} := (c_{i,j_2} - \alpha_{j_1,j_2} - \beta_{j_1,j_2} c_{i,j_1})/\sigma_{j_1,j_2}$, using MAP estimates for instance, where $\alpha_{j_1,j_2}$ and $\sigma_{j_1,j_2}$ are the intercept and scale parameter in the robust model, respectively. One may then flag points with $|z_i^{j_1,j_2}| > 2.5$ (say), which is in line with classical recommendations (see Gervini and Yohai (2002)). Note that the same principle applies for detecting outliers in the columns of $\mathbf{C}$ and, of course, in the multiple linear regressions used afterwards.

Finally note that the percentage of the total variation that is accounted for is capped at 95%, as was the case with traditional PCA. The proposed method may lead to negative eigenvalues as robust correlation matrices are not correlation matrices per se. When this happens, we exclude the associated columns prior to setting $q$.

## 4.2   Robust linear regressions

The convergence result presented at the beginning of Section 4 ensures that posterior model probabilities and estimates of $(\sigma^K, \boldsymbol{\beta}^K)$ based on posterior quantiles (e.g. using posterior medians and Bayesian

credible intervals) are robust to outliers. An analogous convergence result holds for the posterior expectations of the parameters, see Gagnon *et al.* (2018). Predictions for the dependent variable are then obtained by using $\sum_k \pi(k \mid \mathbf{y}) \mathbf{x}_{n+1,k}^T \widehat{\boldsymbol{\beta}}_k$ as in the nonrobust case, the difference being that probabilities and expectations are now computed with respect to the posterior arising from an LPTN error distribution. In Section 4.2.1, we describe the MCMC method used to approximate these probabilities and expectations; in Section 4.2.2, we detail a procedure to efficiently implement this algorithm.

### 4.2.1   Reversible Jump Algorithm

As mentioned in Section 4.1, the price to pay for robustness is an increase in the complexity of the posterior. Parameters are however not restricted to a manifold in the linear regressions. Thus standard numerical approximation methods allow computing integrals with respect to posterior. A commonly employed method for model selection and parameter estimation within the Bayesian paradigm is the RJ algorithm. This sampler allows simulation of the posterior distribution on spaces of varying dimensions, and can thus be used even if the number of parameters in the model is unknown.

The implementation of this sampler requires the specification of some functions, a step typically driven by the structure of the posterior. Recall that, whether there are outliers or not, the posterior under the super heavy-tailed LPTN distribution assumption has a structure similar to that expressed in Proposition 1. In other words, the regression coefficients should be nearly independent given $K$ and $\sigma^K$ and their values should not change dramatically from one model to another. In what follows, we borrow ideas from Gagnon *et al.* (2019), in which an efficient RJ algorithm is built to sample from distributions with similar characteristics.

One iteration of the RJ sampler first randomly selects a model, and then proposes parameters for this model. This candidate model is then accepted as the next state of the Markov chain according to a specific probability; if it is rejected, the chain remains at the same state for another time interval. Specifically, given that the chain currently has $d_K$ components, the sampler that we use randomly selects one of three types of movements: update of the parameters; switch from Model $K$ to Model $K + 1$ (with $d_{K+1} = d_K + 1$); switch from Model $K$ to Model $K - 1$ (with $d_{K-1} = d_K - 1$).

The first step towards obtaining predictions is to identify the statistically relevant PCs. Recall that the individual contribution of each PC is evaluated by comparing the models $I_1 = \{1\}$ and $I_j = \{1, j\}$, $j = 2, \ldots, d$. This first step of the statistical analysis requires $d - 1 = q$ runs of the RJ algorithm that can be performed in parallel. Performing model switches in those RJ samplers thus comes down to adding or withdrawing the $j$-th PC. Denote by $q^*$ the number of PCs associated to Bayes factors greater than the selected threshold; suppose that these statistically significant components are the $j_1$-th, $j_2$-th, $\ldots$, $j_{q^*}$-th PCs. The second step of the analysis then runs a single RJ sampler with $q^* + 1$ nested models, ordered as follows : $I_1 = \{1\}$ (intercept only), $I_2 := \{1, j_1\}, \ldots, I_{q^*+1} := \{1, j_1, \ldots, j_{q^*}\}$. This ensures that the component added (removed) when switching models contains the most (least) information.

The probability mass function used to randomly select the movement type at each iteration is

$$g(j) := \begin{cases} \vartheta, & \text{if } j = 1, \\ (1 - \vartheta)/2, & \text{if } j = 2, 3, \end{cases} \tag{5}$$

where $0 < \vartheta < 1$ is a constant; the value of $\vartheta$ is discussed in Section 4.2.2. At every iteration, an update of the parameters is thus attempted with probability $\vartheta$, while switches to Models $K + 1$ and $K - 1$ are attempted with probability $(1 - \vartheta)/2$ each.

Updating the parameters of Model $K$ is achieved here by using a $(d_K + 1)$-dimensional proposal distribution centred around the current value of the parameter $(\sigma_K, \boldsymbol{\beta}_K)$ and scaled according to $\ell_K$,

where $\ell_K$ is a positive constant given $K$. Each of the $d_K + 1$ candidates is generated independently from the others, according to the one-dimensional strictly positive PDF $\varphi_i, i = 1, \ldots, d_K + 1$. Although the chosen PDF $\varphi_i$ usually is the normal density, we found the PDF in (4) to induce larger candidate steps and to result in a better exploration of the state space. We thus rely on this updating strategy in the analyses of Section 5. Note that one can easily simulate from (4) using the inverse transformation method.

A major issue with the design of RJ algorithms is that there might be a great difference between the "good" values of the parameters under Model $K$ and those under Model $K + 1$ (or $K - 1$). As explained in Section 3, this is not a concern when there is no outlier, or when the same data points are diagnosed as outliers in Models $K$ and $K + 1$; in these cases, the posterior under the LPTN is similar to that under normality. When observations are outliers with respect to Model $K$ but not Model $K + 1$ (say), the posterior of Model $K$ is similar to that under normality excluding outliers, while the posterior of Model $K + 1$ is similar to that under normality based on the whole sample. Therefore, when switching from Model $K$ to Model $K + 1$, the parameters that were already in Model $K$ need to be moved to a position that is appropriate under Model $K+1$. Otherwise, this model switching will be less likely to be accepted, and the sampler will possibly require several iterations before the chain reaches high probability areas. Existing research has focused on that issue and found that it may result in inaccurate estimates, see Brooks *et al.* (2003), Al-Awadhi *et al.* (2004), Hastie (2005), and Karagiannis and Andrieu (2013).

Our strategy for resolving that issue is easily implemented. It consists in adding a vector $\mathbf{c}_{K+1}$ to the current parameters of Model $K$, so as to move these parameters to a suitable area under Model $K + 1$. This leads to a candidate $(\sigma_{K+1}, \boldsymbol{\beta}_{K+1}) := ((\sigma_K, \boldsymbol{\beta}_K) + \mathbf{c}_{K+1}, u_{K+1})$ for Model $K + 1$, where $(\sigma_K, \boldsymbol{\beta}_K)$ is the current value of the parameter under Model $K$ and $u_{K+1}$ is a candidate for the added parameter $\beta_{d_{K+1}, K+1}$, generated from an appropriate strictly positive PDF $q_{K+1}$. To avoid obtaining negative values for $\sigma_K$, we always set the first component of the vectors $\mathbf{c}_i$ to 0.

We now provide a pseudo-code to sample from $\pi(k, \sigma_k, \boldsymbol{\beta}_k \mid \mathbf{y})$ using the RJ sampler. In the next section, we specify the various inputs required to implement this algorithm.

1. Initialise the sampler by setting $(K, \sigma_K, \boldsymbol{\beta}_K)(0)$.
   *Remark*: The number in parentheses beside a vector denotes the iteration.

**Iteration** $m + 1$**.**

2. Generate $u \sim \mathcal{U}(0, 1)$.

   (a) If $u \leq \vartheta$, attempt an update of the parameters. Generate a candidate $\mathbf{w}_{K(m)} := (w_1, \ldots, w_{d_{K(m)}+1})$, where $w_1 \sim \varphi_1(\cdot \mid K(m), \sigma_K(m), \ell_{K(m)})$ and $w_i \sim \varphi_i(\cdot \mid K(m), \beta_{i-1,K}(m), \ell_{K(m)})$ for $i = 2, \ldots, d_{K(m)} + 1$. Generate $u_a \sim \mathcal{U}(0, 1)$; if

   $$u_a \leq \left(1 \wedge \frac{(1/w_1) f(\mathbf{y} \mid K(m), \mathbf{w}_{K(m)})}{(1/\sigma_K(m)) f(\mathbf{y} \mid (K, \sigma_K, \boldsymbol{\beta}_K)(m))}\right),$$

   where

   $$f(\mathbf{y} \mid k, \sigma_k, \boldsymbol{\beta}_k) := \prod_{i=1}^{n} \frac{1}{\sigma_k} f\left(\frac{y_i - \mathbf{x}_{i,k}^T \boldsymbol{\beta}_k}{\sigma_k}\right),$$

   set $(K, \sigma_K, \boldsymbol{\beta}_K)(m + 1) = (K(m), \mathbf{w}_{K(m)})$.

   (b) If $\vartheta < u \leq \vartheta + (1 - \vartheta)/2$, attempt adding a parameter to switch from Model $K(m)$ to Model $K(m) + 1$. Generate $u_{K(m)+1} \sim q_{K(m)+1}$ and $u_a \sim \mathcal{U}(0, 1)$; if

   $$u_a \leq \left(1 \wedge \frac{\pi(K(m) + 1) f(\mathbf{y} \mid K(m) + 1, (\sigma_K, \boldsymbol{\beta}_K)(m) + \mathbf{c}_{K(m)+1}, u_{K(m)+1})}{\pi(K(m)) f(\mathbf{y} \mid (K, \sigma_K, \boldsymbol{\beta}_K)(m)) q_{K(m)+1}(u_{K(m)+1})}\right),$$

set $(K, \sigma_K, \boldsymbol{\beta}_K)(m + 1) = (K(m) + 1, (\sigma_K, \boldsymbol{\beta}_K)(m) + \mathbf{c}_{K(m)+1}, u_{K(m)+1})$.

(c)   If $u > \vartheta + (1 - \vartheta)/2$, attempt withdrawing the last parameter to switch from Model $K(m)$ to Model $K(m) - 1$. Generate $u_a \sim \mathcal{U}(0, 1)$; if

$$u_a \le \left(1 \wedge \frac{\pi(K(m) - 1)f(\mathbf{y} \mid K(m) - 1, (\sigma_K, \boldsymbol{\beta}_{K-})(m) - \mathbf{c}_{K(m)})q_{K(m)}(\beta_{d_K, K}(m))}{\pi(K(m))f(\mathbf{y} \mid (K, \sigma_K, \boldsymbol{\beta}_K)(m))}\right),$$

where $(\sigma_K, \boldsymbol{\beta}_{K-})(m) := (\sigma_K, \beta_{1,K}, \ldots, \beta_{d_{K-1}, K})(m)$, then set $(K, \sigma_K, \boldsymbol{\beta}_K)(m+1) = (K(m)-1, (\sigma_K, \boldsymbol{\beta}_{K-})(m) - \mathbf{c}_{K(m)})$.

3.   In case of rejection, set $(K, \sigma_K, \boldsymbol{\beta}_K)(m + 1) = (K, \sigma_K, \boldsymbol{\beta}_K)(m)$.

4.   Go to Step 2.

It is easily verified that the resulting stochastic process $\{(K, \sigma_K, \boldsymbol{\beta}_K)(m) : m \in \mathbb{N}\}$ is a $\pi(k, \sigma_k, \boldsymbol{\beta}_k \mid \mathbf{y})$-irreducible and aperiodic Markov chain. Furthermore, it satisfies the reversibility condition with respect to the posterior, as stated in the following proposition. Therefore, it is an ergodic Markov chain, which guarantees that the Law of Large Numbers holds.

**Proposition 2.** *The Markov chain* $\{(K, \sigma_K, \boldsymbol{\beta}_K)(m) : m \in \mathbb{N}\}$ *arising from the RJ described above satisfies the reversibility condition with respect to the posterior* $\pi(k, \sigma_k, \boldsymbol{\beta}_k \mid \mathbf{y})$.

*Proof.* See the supplementary material (Section 7).                                      □

### 4.2.2   Efficient implementation

An optimal implementation of the RJ algorithm described above requires carefully selecting the various inputs: the PDFs $q_i$, the constants $\vartheta$ and $\ell_i$, and the vectors $\mathbf{c}_i$. Hereafter, "optimal implementation" or "optimal design" means that the generated Markov chain mixes as rapidly as possible, thus engendering least variable estimators.

In Gagnon *et al.* (2019), a posterior structure similar to that expressed in Proposition 1 is considered, and theoretical results leading to an optimal RJ algorithm are obtained. In that paper, the parameters of any given model are conditionally independent and identically distributed. An implicit assumption on the posterior studied is that distributions of parameters remain the same when switching from Model $K$ to Model $K + 1$ (or $K - 1$). The authors find asymptotically optimal values for $\vartheta$ and $\ell_K$ (as the number of parameters approaches infinity). They conjecture that their results are valid (to some extent) when the parameters are conditionally independent, but not identically distributed (for any given model). They also provide guidelines to suitably design the PDFs $q_K$. We use these results as a starting point in the design of our RJ algorithm.

In the settings of Gagnon *et al.* (2019), the asymptotically optimal value for $\vartheta$ depends on the PDFs $q_i$. It is also empirically observed that for moderate values of $K_{\max}$, selecting any value between 0.2 and 0.6 is almost optimal. We use $\vartheta := 0.6$ in the numerical analyses of Section 5, as $K_{\max}$ is rather small (there are few models to visit). Generally speaking, larger values of $\vartheta$ leave the chain more time for exploring the parameters' state space between model switches. Based on several runs of the RJ algorithm, $\vartheta := 0.6$ is in fact nearly optimal for the data in Section 5.

If the parameters $(\sigma_K, \boldsymbol{\beta}_K)$ were independent and identically distributed for each model, the asymptotically optimal value for $\ell_K$ would be $\ell/\sqrt{d_K + 1}$, with $\ell$ tuned to accept 23.4% of candidates $\mathbf{w}_K$. When these assumptions are violated, the asymptotically optimal value for $\ell$ usually corresponds to an acceptance rate smaller than 0.234 (see Bédard (2007) and Bédard (2019)). Considering this, and

adding the fact that $d_k$ may be rather small, we recommend to perform trial runs to identify optimal values for all $\ell_k$. We use the 0.234 rule within each model to initiate the process. In our analyses in Section 5, the optimal values for all $\ell_k$ correspond to an acceptance rate relatively close to 0.234.

We propose to specify the PDFs $q_i$ and vectors $\mathbf{c}_i$ through trial runs as well. Specifying these functions and vectors requires information about locations and scalings of regression coefficients for all models. We gather this information by running a random walk Metropolis algorithm for each model; this sampler may be seen as a RJ algorithm in which $\vartheta := 1$ (i.e. a sampler in which only updates of the parameters are proposed). The recommended procedure is now detailed.

**For each $k \in \{1, \ldots, K_{\max}\}$:**

1. Tune $\ell_k$ such that the acceptance rate of candidates $\mathbf{w}_k$ is approximately 0.234; denote this value by $\ell_k^{\text{start}}$.

2. Select a sequence of values around $\ell_k^{\text{start}}$: $(\ell_{1,k}, \ldots, \ell_{j_0,k} := \ell_k^{\text{start}}, \ldots, \ell_{L,k})$, where $L$ is a positive integer.

3. For each $\ell_{j,k}$, run a random walk Metropolis sampler initialised as follows: $(\sigma_k(0))^2 \sim \text{Inv-}\Gamma$ with shape and rate given by $(n - d_k)/2$ and $\|\mathbf{y} - \widehat{\mathbf{y}}_k\|_2^2/2$, respectively; $\beta_{1,k}(0) \sim \mathcal{N}(\widehat{\beta}_{1,k}, (\sigma_k(0))^2/n)$, and $\beta_{j,k}(0) \sim \mathcal{N}(\widehat{\beta}_{j,k}, (\sigma_k(0))^2/(n-1))$, $j = 2, \ldots, d_k$ (if $k \geq 2$). Here, $\widehat{\mathbf{y}}_k$ is computed using a preliminary robust estimate $\widehat{\boldsymbol{\beta}} := (\widehat{\beta}_{1,k}, \ldots, \widehat{\beta}_{d_k,k})$ (MAP estimate under the robust LPTN model for instance).

4. For each $\ell_{j,k}$ ($j = 1, \ldots, L$), estimate the location and scaling of each $\beta_{i,k}$ using the runs in Step (3). In particular, compute the mean (denoted by $m_{i,j}^k$) and standard deviation (denoted by $s_{i,j}^k$) of $\{\beta_{i,k}(m) : m \in \{B + 1, \ldots, T\}\}$, for $i = 1, \ldots, d_k$, where $B$ is the length of the burn-in period and $T$ the number of iterations. Repeat for $\sigma_k$, denoting the means and standard deviations by $m_{\sigma,j}^k$ and $s_{\sigma,j}^k$. Measure the efficiency of the sampler with respect to $\ell_{j,k}$ using the sum of the integrated autocorrelation times (IAT) of $\{\sigma_k(m) : m \in \{B + 1, \ldots, T\}\}$ and $\{\beta_{i,k}(m) : m \in \{B + 1, \ldots, T\}\}$ for $i = 1, \ldots, d_k$. Record the value $\ell_k^{\text{opt}}$ corresponding to the smallest IAT.

5. If $\ell_k^{\text{opt}}$ corresponds to the lower or upper bound of the range defined in Step (2), i.e. $\ell_{1,k}$ or $\ell_{L,k}$, change the sequence of values for $\ell_k$ and repeat.

6. For $i = 1, \ldots, d_k$, compute the average of $\{m_{i,1}^k, \ldots, m_{i,L}^k\}$ (denoted by $m_{i,k}$) and $\{s_{i,1}^k, \ldots, s_{i,L}^k\}$ (denoted by $s_{i,k}$). Also compute the average of $\{m_{\sigma,1}^k, \ldots, m_{\sigma,L}^k\}$ (denoted by $m_{\sigma,k}$) and $\{s_{\sigma,1}^k, \ldots, s_{\sigma,L}^k\}$ (denoted by $s_{\sigma,k}$).

These runs can be performed in parallel for computational efficiency, in an automatic procedure that allows users to retrieve the desired output at the end. Using this output, set $q_j$ ($j = 2, \ldots, K_{\max}$) equal to the distribution in (4), with location and scale parameters given by $m_{d_j,j}$ and $s_{d_j,j}$, respectively. Also set $\mathbf{c}_j := (0, m_{1,j} - m_{1,j-1}, \ldots, m_{d_j-1,j} - m_{d_j-1,j-1})^T$, $j = 2, \ldots, K_{\max}$, and $\ell_k$ equal to $\ell_k^{\text{opt}}$ for all $k$.

The only inputs left to choose before implementing the RJ algorithm are the initial values for the model indicator and parameters. We recommend to generate $K(0) \sim \mathcal{U}\{1, \ldots, K_{\max}\}$, $\sigma_K(0)$ from a normal truncated at 0 with mean $m_{\sigma,K(0)}$ and standard deviation $s_{\sigma,K(0)}$, and $\beta_{j,K(0)}(0) \sim \mathcal{N}(m_{j,K(0)}, s_{j,K(0)}^2)$, $j = 1, \ldots, d_{K(0)}$. The analyses of Section 5 rely on sequences of length $L = 11$ for $\ell_k$ ($\ell_k^{\text{start}}$ is the median), $T = 100{,}000$ iterations, and a burn-in period of length $B = 10{,}000$ for the trial runs. When running the RJ sampler, we use 1,000,000 iterations and a burn-in period of length 100,000.

# 5    Case study: prediction of returns for the S&P 500

In this section, we illustrate the performance of our robust approach on a real data set containing outliers. We provide a detailed analysis of the results and contrast them with those from other approaches to identify in which situations it is expected to perform better. The data set and context are described in Section 5.1, the competitors are presented in Section 5.2, while the section finishes with the result analysis and comparison in Section 5.3.

The use of super heavy-tailed distributions in linear regression has recently been introduced in Desgagné and Gagnon (2019), where the special case of simple linear regressions through the origin was studied. The usual linear regression model was later analysed in Gagnon *et al.* (2018). Although theoretical results about model selection are presented in Gagnon *et al.* (2018), it is the first time that an illustration of the practical benefits is presented in that context.

## 5.1    Data set and context description

In this example, we model the January 2011 daily returns of the S&P 500 by exploiting their potential linear relationship with some financial assets and indicators. We next use the estimated models to predict the February 2011 daily returns of this stock index; to this end, covariate observations on day $i$ will be used to predict the return of the S&P 500 on day $i + 1$. A detailed list of the 18 covariates considered is provided in the supplementary material (Section 7); $n = 19$ observations are available for model estimation. The full linear regression model with all covariates would have 20 parameters ($p = 18$ regression coefficients for the covariates, to which we add the intercept and scale parameter). We perform robust and nonrobust PCA procedures, which are expected to be beneficial given that financial assets and indicators are likely to carry redundant information.

## 5.2    Competitors

The results are compared with those obtained under the normality of errors assumption (nonrobust Bayesian approach) to evaluate outlier protection performance. The classical frequentist approach and the robust frequentist approach of Hubert and Verboven (2003) are also included in the comparison. The implementation of a robust frequentist approach allows contrasting the effects of our Bayesian robust PCA decomposition. The proposed model-based PC selection approach is also evaluated.

In principle, to construct a PCR approach, one only needs a PCA and a linear regression method. There are of course many combinations of PCA and linear regression approaches possible. To keep the analysis and comparison simple, we restrict our attention to a single combination of these methods for each of the four classes of PCR approaches considered (robust/nonrobust and Bayesian/frequentist). The four combinations selected are arguably the best approaches in each of the four classes. In the robust frequentist approach, we use MM-regression (Yohai, 1987) as it offers one of the best available asymptotic *breakdown point versus efficiency* tradeoffs. We nevertheless acknowledge that there exist several other good combinations; one could, for instance, use least trimmed squares estimators (LTS, Rousseeuw (1985)), M-estimators (Huber, 1973), S-estimators (Rousseeuw and Yohai, 1984), or other more recent robust regression approaches like those of Agostinelli and Greco (2013) and Atkinson *et al.* (2017). We refer the reader to the recent robust regression comparisons presented in Gagnon *et al.* (2018) and Yu and Yao (2017), which respectively focus on Bayesian and frequentist methods.

Model estimation is performed using each of the four mentioned approaches. In the robust and nonrobust Bayesian approaches, statistically significant PCs are identified by relying on Bayes factors

with a threshold of 1. This means that when on average (over the parameter space) an individual PC improves the fit over the model consisting solely of the intercept, then it is included in the second stage of the statistical analysis (for building the nested models).
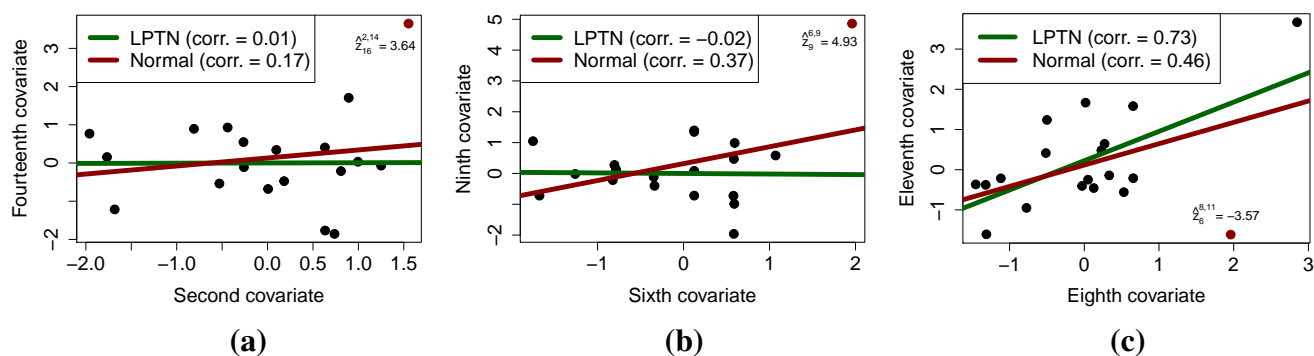
## 5.3 Results and analysis

The average absolute deviations (AAD) between the predicted and actual February 2011 returns are reported in the first column of Table 1. In the current financial context, it may also be of interest to predict whether the asset (S&P 500 in our case) will go up or down the following day. Using the sign of our predicted returns, we find to be correct 10, 13, 10, and 11 times out of 19 under the normal, LPTN, classical and robust frequentist approaches; the success rates are reported in the second column of Table 1.

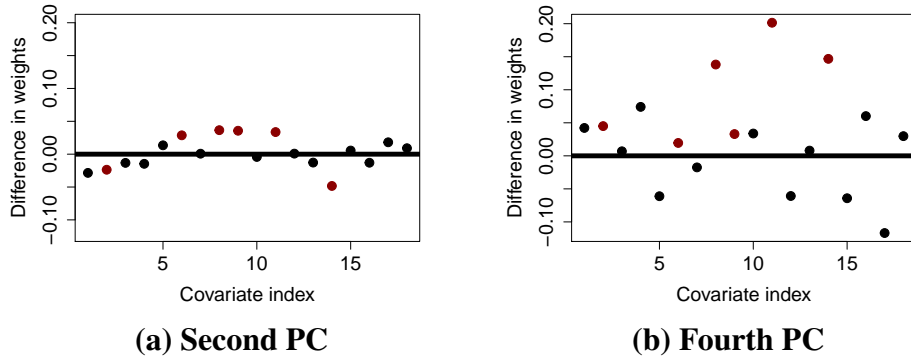| Approach | AAD | Sign prediction rate |
|---|---|---|
| Normal errors (nonrobust Bayesian) | 0.60 | 0.53 |
| LPTN errors (robust Bayesian) | 0.49 | 0.68 |
| Classical frequentist | 0.63 | 0.53 |
| Robust frequentist | 0.57 | 0.58 |

**Table 1.** Prediction results for the February 2011 daily returns of the S&P 500, using robust and nonrobust versions of Bayesian and frequentist approaches

We know that differences in the results obtained from the normal and LPTN models are essentially due to the presence of outliers. The outlier detection method described in Section 4.1 indeed flags several observations, both in the PCA and linear regression steps. Each graph in Figure 3 illustrates linear relationships between a different pair of covariates; the two linear relationships in a given graph are established using the normal and LPTN error distributions. We see that in the presence of outliers, the choice of error distribution obviously has a large impact on the trends obtained from the data. Figure 5 also depicts linear relationships using the normal and LPTN distributions, but this time between the dependent variable and some PCs. Specifically, the graphs on the top line picture linear relationships between the dependent variable and the second PC; in the left graph, the second PC was constructed using a robust PCA while in the right graph, that same PC was constructed using a traditional PCA. The exercise is then repeated with the fourth PC and produces the two graphs on the bottom line.



**Figure 3.** Linear relationships between the (a) second and fourteenth covariates; (b) sixth and ninth covariates; (c) eighth and eleventh covariates

From these figures, it is clear that contaminated correlation estimation in traditional PCA leads to an inferior assessment of the relationships between covariates (Figure 3), which in turn leads to a different way of constructing the PCs (Figure 4). By *inferior*, we mean here that the trend does not reflect the behaviour of the majority of the observations, but rather consists in a poor compromise between that behaviour and the behaviour of outliers. The left graph of Figure 4 plots the differences arising from applying a robust PCA rather than a traditional one when computing the covariates weights used to construct the second PC. The right graph repeats the exercise for the fourth PC, which leads to even greater differences than for the second PC.



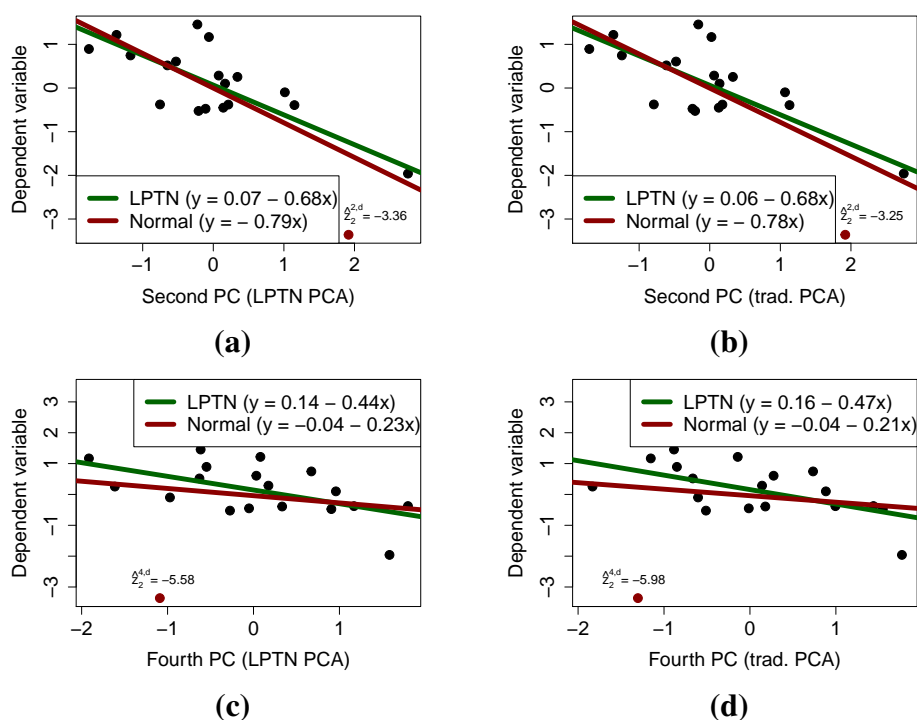**(a) Second PC**                    **(b) Fourth PC**

**Figure 4.** Differences in covariate weights used to construct the (a) second PC and (b) fourth PC, when comparing the robust and traditional PCA; the six covariates considered in Figure 3 are shown in red (2nd, 6th, 8th, 9th, 11th, and 14th covariates)

These discrepancies are ultimately seen to have an impact on the estimated linear regressions (Figure 5). In particular, combining the traditional PCA (instead of the robust one) with LPTN regression models gives an average absolute deviation of 0.53 (instead of the value 0.49 appearing in Table 1). When implementing the robust PCR, only the second and fourth PC end up being retained for the second stage of the statistical analysis. This means that the models corresponding to $I_1 := \{1\}$ (intercept only), $I_2 := \{1, 2\}$, and $I_3 := \{1, 2, 4\}$ are the only models considered in that stage. Their posterior probabilities each are 0.00, 0.18, and 0.82. Note that the fourth PC is not selected by the normal (nonrobust) approach.

The difference between Bayesian and frequentist robust approaches not only resides in the selection of statistically significant PCs (subsequently used in prediction), but also in the PCA decomposition employed. In fact, if the robust frequentist approach were to use the same regressors as the Bayesian one, and then retaining the usual approach for selecting PCs and estimating parameters, the results would be (essentially) the same as those of the robust Bayesian approach. The results would also be (essentially) the same if the exact same PCs as those in the Bayesian approach were selected (the frequentist approach selects an additional one, as explained in the next paragraph). Our analysis thus shows that, in our example, the proposed PCA represents the information from the covariates in lower dimensional spaces in a way that is more suitable to predict the dependent variable.

Our analysis also shows that, in the current example, the Bayesian approach for selecting statistically significant PCs does not dominate the frequentist one, and vice versa. The frequentist approach however leads to an overparameterised model, an undesirable characteristic. Indeed, frequentist approaches use cross-validation, along with a robust prediction measure in the case of the robust method. Models with an increasing number of PCs are thus evaluated (the PCs are ordered, and included in the models according to this predetermined order); the model enjoying the best fit is then used for prediction. If we apply this method using the PCs obtained from our robust PCA for instance, we find the

**Figure 5.** Linear relationships between the dependent variable and second PC constructed using the (a) robust PCA and (b) traditional PCA; linear relationships between the dependent variable and fourth PC constructed using the (c) robust PCA and (d) traditional PCA

model with four PCs to be the best option, and so according to this analysis, the first four PCs should be used for prediction. If we consider a larger class of models instead of being forced to include PCs in a predetermined manner, we however find through cross-validation methods that the model with the first, second and fourth PCs performs better (as pointed out by the proposed Bayesian method). Indeed, the third PC does not significantly explain the variability in the dependent variable; it may thus negatively influence the model's ability to generalise (see Hadi and Ling (1998) and Jolliffe (1982) for examples).

We note that the frequentist approach could be modified so as to include PCs on the merit of their individual contribution. To this effect, the Bayesian information criterion (BIC) could be used to evaluate the individual contributions of the PCs, after which the frequentist method could be applied on the retained ordered components only; the resulting set of PCs would be the one used in prediction. Similarly, our robust Bayesian PCR could also be applied under the frequentist paradigm.

We also ran simulations and drew the same conclusions as in this section. We thus do not present them for brevity. We only note that the robust approaches (Bayesian and frequentist) are expected to be efficient by their nature, and this is what we observed; they perform only slightly worse than their nonrobust counterparts when there are no outliers.

# 6   Conclusion and further remarks

In light of the results of Section 5, we conclude that the proposed robust Bayesian PCR approach is expected to perform better than its competitors (at least those that are nonrobust) when there are outliers in the data set (either among the covariates or the dependent variable). This is a consequence of the new class of super heavy-tailed PCA models, combined to the LPTN regressions of Gagnon *et al.*

(2018). The approach is also expected to perform better when the first $q$ PCs do not all contribute in explaining the variability of the dependent variable. The approach indeed takes advantage of the linear independence of the PCs to effectively exclude the components that are not relevant, and next forms a sequence of nested models from which predictions are produced and averaged out according to the model posterior probabilities.

As explained in Section 4.1, the robust PCA applied to the real data of Section 5 is, in reality, an approximation to the exact wholly robust PCA model. Further research is needed to acquire a deeper understanding of its theoretical properties, as well as to develop an efficient implementation method. It would be particularly useful to obtain a robust procedure that not only reduces dimensionality, but also induces sparsity to deal with cases where $p \gg n$.

# Acknowledgements

# References

Agostinelli, C. and Greco, L. (2013) A weighted strategy to handle likelihood uncertainty in Bayesian inference. *Comput. Statist.*, **28**, 319–339.

Al-Awadhi, F., Hurn, M. and Jennison, C. (2004) Improving the acceptance rate of reversible jump MCMC proposals. *Statist. Probab. Lett.*, **69**, 189–198.

Andrade, J. A. A. and O'Hagan, A. (2011) Bayesian robustness modeling of location and scale parameters. *Scand. J. Stat.*, **38**, 691–711.

Atkinson, A. C., Corbellini, A. and Riani, M. (2017) Robust Bayesian regression with the forward search: theory and data analysis. *TEST*, **26**, 869–886.

Bédard, M. (2007) Weak convergence of Metropolis algorithms for non-i.i.d. target distributions. *Ann. Appl. Probab.*, **17**, 1222–1244.

— (2019) Hierarchical models and tuning of random walk Metropolis algorithms. *Journal of Probability and Statistics*, 1–24.

Brooks, S. P., Giudici, P. and Roberts, G. O. (2003) Efficient construction of reversible jump Markov chain Monte Carlo proposal distributions. *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, **65**, 3–39.

Casella, G., Giròn, F. J., Martínez, M. L. and Moreno, E. (2009) Consistency of Bayesian procedures for variable selection. *Ann. Statist.*, **37**, 1207–1228.

Desgagné, A. (2015) Robustness to outliers in location–scale parameter model using log-regularly varying distributions. *Ann. Statist.*, **43**, 1568–1595.

Desgagné, A. and Gagnon, P. (2019) Bayesian robustness to outliers in linear regression and ratio estimation. *Braz. J. Probab. Stat.*, **33**, 205–221. ArXiv:1612.05307.

Gagnon, P., Bédard, M. and Desgagné, A. (2019) Weak convergence and optimal tuning of the reversible jump algorithm. *Math. Comput. Simulation*, **161**, 32–51.

Gagnon, P., Desgagné, A. and Bédard, M. (2018) A new Bayesian approach to robustness against outliers in linear regression. *Bayesian Anal.* Advance publication.

George, E. I. and McCulloch, R. E. (1993) Variable selection via Gibbs sampling. *J. Amer. Statist. Assoc.*, **88**, 881–889.

Gervini, D. and Yohai, V. J. (2002) A class of robust and fully efficient regression estimators. *Ann. Statist.*, **30**, 583–616.

Green, P. J. (1995) Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika*, **82**, 711–732.

Hadi, A. S. and Ling, R. F. (1998) Some cautionary notes on the use of principal components regression. *Amer. Statist.*, **52**, 15–19.

Hastie, D. (2005) *Towards Automatic Reversible Jump Markov Chain Monte Carlo*. Ph.D. thesis, University of Bristol.

Hoeting, J. A., Madigan, D., Raftery, A. E. and Volinsky, C. T. (1999) Bayesian model averaging: A tutorial. *Statist. Sci.*, 382–401.

Huber, P. J. (1973) Robust regression: Asymptotics, conjectures and monte carlo. *Ann. Statist.*, 799–821.

Hubert, M. and Verboven, S. (2003) A robust PCR method for high-dimensional regressors. *Journal of Chemometrics: A Journal of the Chemometrics Society*, **17**, 438–452.

Jeffreys, H. (1967) *Theory of Probability*. Oxford Univ. Press, London.

Jolliffe, I. (2011) *Principal component analysis*. Springer.

Jolliffe, I. T. (1982) A note on the use of principal components in regression. *J. R. Stat. Soc. Ser. C. Appl. Stat.*, **31**, 300–303.

Karagiannis, G. and Andrieu, C. (2013) Annealed importance sampling reversible jump MCMC algorithms. *J. Comp. Graph. Stat.*, **22**, 623–648.

Lindley, D. V. (1957) A statistical paradox. *Biometrika*, **44**, 187–192.

Luttinen, J., Ilin, A. and Karhunen, J. (2009) Bayesian robust PCA for incomplete data. In *International conference on independent component analysis and signal separation*, 66–73.

Raftery, A. E., Madigan, D. and Hoeting, J. A. (1997) Bayesian model averaging for linear regression models. *J. Amer. Statist. Assoc.*, **92**, 179–191.

Rousseeuw, P. J. (1985) Multivariate estimation with high breakdown point. *Mathematical statistics and applications*, **37**, 283–297.

Rousseeuw, P. J. and Yohai, V. J. (1984) Robust regression by means of s-estimators. In *Robust and Nonlinear Time Series Analysis.*, 256–272. Springer.

Schwarz, G. (1978) Estimating the dimension of a model. *Ann. Statist.*, **6**, 461–464.

Šmídl, V. and Quinn, A. (2007) On Bayesian principal component analysis. *Comput. Statist. Data Anal.*, **51**, 4101–4123.

Tipping, M. E. and Bishop, C. M. (1999) Probabilistic principal component analysis. *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, **61**.

Tipton, J., Hooten, M. and Goring, S. (2017) Reconstruction of spatio-temporal temperature from sparse historical records using robust probabilistic principal component regression. *Adv. Stat. Clim. Meteorol. Oceanogr.*, **3**, 1–16.

Wang, L. (2012) Bayesian principal component regression with data-driven component selection. *J. Appl. Stat.*, **39**, 1177–1189.

West, M. (1984) Outlier models and prior distributions in Bayesian linear regression. *J. R. Stat. Soc. Ser. B. Stat. Methodol.*, **46**, 431–439.

— (2003) Bayesian factor regression models in the "large p, small n" paradigm. In *Bayesian Statistics 7*, 723–732. Oxford Univ. Press, London.

Yohai, V. J. (1987) High breakdown-point and high efficiency robust estimates for regression. *Ann. Statist.*, **15**, 642–656.

Yu, C. and Yao, W. (2017) Robust linear regression: A review and comparison. *Comm. Statist. B – Simulation Comput.*, **46**, 6261–6282.

Zhao, Q., Meng, D., Xu, Z., Zuo, W. and Zhang, L. (2014) Robust principal component analysis with complex noise. In *International conference on machine learning*, 55–63.

# 7    Supplementary material

We first present the mathematical justification of the approximate robust principal component analysis (PCA) in Section 7.1. The validity of our prior structure is next addressed in Section 7.2. Propositions 1 and 2 are proved in Section 7.3. The list of the explanatory variables considered in the real data analysis in Section 5 of our paper is provided in Section 7.4.

## 7.1 Mathematical justification of the approximate robust PCA

See Section 4.1 of our paper for the definition of notation. Given that the LPTN matches the normal distribution everywhere except in the tails, the limiting posterior of $(\tilde{\mathbf{Z}}_q, \tilde{\mathbf{L}}_q, \tilde{\mathbf{A}}_q, \eta)$ based on the exact robust PCA is similar to that arising from the traditional PCA model with normal errors based on $\mathbf{C}^*$, as the outliers moves away from the general trend. This means that the exact robust PCA applied to $\mathbf{C}$ leads to essentially the same singular value decomposition as a traditional PCA applied to $\mathbf{C}^*$ (in the limit). The approximate robust PCA method relies on this equivalence.

The first step in performing an approximate robust PCA is to obtain $\mathbf{C}$ by standardising the columns of the original data set. Location and scale estimates $\widehat{\mu}_j$ and $\widehat{\sigma}_j$ are thus used to standardise Column $j$, $j = 1, \ldots, p$. Relying on a robust location-scale model as in Desgagné (2015), with an LPTN error distribution and $\rho := 0.95$, ensures that $(\widehat{\mu}_j, \widehat{\sigma}_j) \longrightarrow (\widehat{\mu}_j^{-O}, \widehat{\sigma}_j^{-O})$, where $(\widehat{\mu}_j^{-O}, \widehat{\sigma}_j^{-O})$ are estimates based on nonoutliers only. Note that the robust location-scale model is the linear regression model with only the intercept. For large $n$, we also have $(\widehat{\mu}_j^{-O}, \widehat{\sigma}_j^{-O}) \approx (\widehat{\mu}_j^*, \widehat{\sigma}_j^*)$, where $(\widehat{\mu}_j^*, \widehat{\sigma}_j^*)$ are the sample mean and standard deviation obtained from $\mathbf{C}^*$, which is based on the normality of errors and in which outliers are replaced by their vertical projection. Denote by $c_{ij}^O$ the outlying values; they are excluded for the estimation of $(\widehat{\mu}_j^{-O}, \widehat{\sigma}_j^{-O})$ and replaced by their vertical projection, denoted by $c_{ij}^*$, for the estimation of $(\widehat{\mu}_j^*, \widehat{\sigma}_j^*)$. Provided that $n$ is large enough, the impact of those points on the sample mean and standard deviation will indeed be negligible; furthermore, it was previously argued that estimates obtained under LPTN and normal error distributions are similar. The resulting matrices $\mathbf{C}$ and $\mathbf{C}^*$ are the same in the limit, except on lines containing outliers.

The second step in performing the approximate robust PCA consists in computing robust correlations between all pairs of columns in $\mathbf{C}$. We know that the correlation between the standardised columns $j_1$ and $j_2$ of $\mathbf{C}^*$ is $\widehat{\beta}_{j_1,j_2}^N$, the OLS slope estimate. We are interested in comparing the robust slope estimator (applied to columns of the matrix $\mathbf{C}$) to $\widehat{\beta}_{j_1,j_2}^N$. When using a robust regression model as in Gagnon *et al.* (2018) with an LPTN error distribution and $\rho := 0.95$, we find $\widehat{\beta}_{j_1,j_2} \longrightarrow \widehat{\beta}_{j_1,j_2}^{-O}$, where $\widehat{\beta}_{j_1,j_2}^{-O}$ is the robust slope estimate obtained using nonoutliers only. Again, for large $n$, we find $\widehat{\beta}_{j_1,j_2}^{-O} \approx \widehat{\beta}_{j_1,j_2}^N$. The robust correlation matrix obtained from $\mathbf{C}$ is thus asymptotically equal to the correlation matrix obtained from $\mathbf{C}^*$. Its diagonal elements are equal to 1; for simplicity, we set the upper diagonal entries to $\widehat{\beta}_{j_1,j_2}$, where Column $j_2$ plays the role of the dependent variable; we then make the matrix symmetrical.

The PCs $\widehat{\mathbf{Z}}_q$ are ultimately computed using $\mathbf{C}\widehat{\mathbf{v}}_j$, with $\widehat{\mathbf{v}}_j$ being the $j$-th eigenvector of the robust correlation matrix of $\mathbf{C}$.

## 7.2 Validity of our prior structure

Relying on improper priors such as $\pi(\sigma_k, \boldsymbol{\beta}_k \mid k) = c_k/\sigma_k$ may lead to inconsistencies in model selection (see Casella *et al.* (2009)). For instance, one could select different constants $c_k$ in different models so as to yield the desired conclusions. In this section, we show that the Jeffreys-Lindley paradox does not arise in our PCR context under the normal distribution assumption. It is thus expected to not arise either under the robust LPTN distribution, given its similarity to the normal.

Consider Models $s$ and $t$, where Model $s$ is nested in Model $t$. The ratio of the posterior probabilities of these two models is given by (see Proposition 3.1 in our paper)

$$\frac{\pi(t \mid \mathbf{y})}{\pi(s \mid \mathbf{y})} = \frac{\Gamma((n - d_s)/2 - (d_t - d_s)/2)}{\Gamma((n - d_s)/2)((n - d_s)/2)^{-(d_t - d_s)/2}} \, n^{-(d_t - d_s)/2} \left( \frac{\|\mathbf{y} - \widehat{\mathbf{y}}_s\|^2/(n - 1)}{\|\mathbf{y} - \widehat{\mathbf{y}}_t\|^2/(n - 1)} \right)^{n/2}$$

$$\times \frac{\pi^{d_t/2}}{\pi^{d_s/2}} \frac{((n-d_s)/2)^{-(d_t-d_s)/2}}{n^{-(d_t-d_s)/2}} \frac{\left(\|\mathbf{y} - \widehat{\mathbf{y}}_t\|^2/(n-1)\right)^{d_t/2}}{\left(\|\mathbf{y} - \widehat{\mathbf{y}}_s\|^2/(n-1)\right)^{d_s/2}} \frac{\pi(t)}{\pi(s)}. \tag{6}$$

The difference between the Bayesian information criteria (BIC, Schwarz (1978)) of Models $t$ and $s$ is given by

$$\begin{aligned}
\text{BIC}_t - \text{BIC}_s &= n \log\left(\|\mathbf{y} - \widehat{\mathbf{y}}_t\|^2/n\right) + (d_t + 1)\log n \\
&\quad - n \log\left(\|\mathbf{y} - \widehat{\mathbf{y}}_s\|^2/n\right) - (d_s + 1)\log n \\
&= n \log\left(\frac{\|\mathbf{y} - \widehat{\mathbf{y}}_t\|^2/n}{\|\mathbf{y} - \widehat{\mathbf{y}}_s\|^2/n}\right) + (d_t - d_s)\log n.
\end{aligned}$$

Given that the first ratio on the right hand side of (6) converges to 1 as $n \longrightarrow \infty$, we have that $\exp\{-(\text{BIC}_t - \text{BIC}_s)/2\}$ asymptotically behaves like the first part on the right hand side of (6). The terms $\left(\|\mathbf{y} - \widehat{\mathbf{y}}_k\|^2/(n-1)\right)^{d_k/2}$ in (6) converge towards a constant (in $n$) and are thus dominated. The other terms in (6) are either constant in terms of $n$ or dominated as well. Therefore, $\pi(t|\mathbf{y})/\pi(s|\mathbf{y})$ and $\exp\{-(\text{BIC}_t - \text{BIC}_s)/2\}$ share the same asymptotic behaviour. This will be sufficient to prove that the prior structure does not prevent the Bayesian variable selection procedure to be consistent, in the same sense as Casella *et al.* (2009). If the "true" model is among the models considered, then its posterior probability converges to 1 as $n$ increases. Further technical details are required for a rigorous proof. Empirical evidences also point towards the validity of our claim.

It would be interesting to investigate the asymptotic behaviour in the more general context of traditional linear regression. The fact that the regressors are standardised and linearly independent plays a role in the sketch of the proof presented above. It would however be surprising if a similar prior structure, but with slightly correlated standardised regressors, led to inconsistencies.

In practice (with finite samples), one may set the prior $\pi(k)$ to be proportional to $\pi^{-d_k/2}$ times a prior opinion about $\left(\|\mathbf{y} - \widehat{\mathbf{y}}_k\|^2/(n-1)\right)^{-d_k/2}$, to cancel the effect of these two terms in (6). In the numerical analyses, we set $\pi(k) \propto 1$ because we do not have relevant information. Note that the robust approach proposed in this paper can be used with any informative prior such as those in Raftery *et al.* (1997).

## 7.3 Proofs

*Proof of Proposition 3.1.* The proof is essentially a computation using that $f := \mathcal{N}(0,1)$ and the structure of the principal components. First,

$$\begin{aligned}
\pi(k, \sigma_k, \boldsymbol{\beta}_k \mid \mathbf{y}) &\propto f(\mathbf{y} \mid k, \sigma_k, \boldsymbol{\beta}_k)\pi(\sigma_k, \boldsymbol{\beta}_k \mid k)\pi(k) \\
&\propto f(\mathbf{y} \mid k, \sigma_k, \boldsymbol{\beta}_k)(1/\sigma_k)\pi(k).
\end{aligned}$$

The likelihood function for a given model is

$$\begin{aligned}
f(\mathbf{y} \mid k, \sigma_k, \boldsymbol{\beta}_k) &= \prod_{i=1}^{n} \frac{1}{\sigma_k \sqrt{2\pi}} \exp\left\{-\frac{1}{2\sigma_k^2}(y_i - \mathbf{x}_{i,k}^T \boldsymbol{\beta}_k)^2\right\} \\
&= \frac{1}{\sigma_k^n (2\pi)^{n/2}} \exp\left\{-\frac{1}{2\sigma_k^2} \sum_{i=1}^{n}(y_i - \mathbf{x}_{i,k}^T \boldsymbol{\beta}_k)^2\right\}.
\end{aligned}$$

We now analyse the sum in the exponential:

$$\sum_{i=1}^{n}(y_i - \mathbf{x}_{i,k}^T\boldsymbol{\beta}_k)^2 = \sum_{i=1}^{n}y_i^2 - 2\sum_{i=1}^{n}y_i\sum_{j=1}^{d_k}x_{iI_{j,k}}\beta_{j,k} + \sum_{i=1}^{n}\left(\sum_{j=1}^{d_k}x_{iI_{j,k}}\beta_{j,k}\right)^2$$

$$= n - 1 - 2\sum_{j=1}^{d_k}\beta_{j,k}\sum_{i=1}^{n}y_i x_{iI_{j,k}} + \sum_{i=1}^{n}\left(\sum_{j=1}^{d_k}x_{iI_{j,k}}\beta_{j,k}\right)^2,$$

using that $\sum_{i=1}^{n}y_i^2 = n - 1$. We also have

$$\sum_{i=1}^{n}\left(\sum_{j=1}^{d_k}x_{iI_{j,k}}\beta_{j,k}\right)^2 = \sum_{i=1}^{n}\left(\sum_{j=1}^{d_k}(x_{iI_{j,k}}\beta_{j,k})^2 + \sum_{j,s=1(j\neq s)}^{d_k}x_{iI_{j,k}}\beta_{j,k}x_{iI_{s,k}}\beta_{s,k}\right)$$

$$= \sum_{j=1}^{d_k}\beta_{j,k}^2\sum_{i=1}^{n}x_{iI_{j,k}}^2,$$

using $\sum_{i=1}^{n}x_{ij}x_{is} = 0$ for all $j, s \in \{2, \ldots, d\}$ with $j \neq s$, $x_{11} = \ldots = x_{n1} = 1$, $(1/n)\sum_{i=1}^{n}x_{ij} = 0$ for all $j \in \{2, \ldots, d\}$. Consequently,

$$\sum_{i=1}^{n}(y_i - \mathbf{x}_{i,k}^T\boldsymbol{\beta}_k)^2 = n - 1 - 2\sum_{j=1}^{d_k}\beta_{j,k}\sum_{i=1}^{n}y_i x_{iI_{j,k}} + \sum_{j=1}^{d_k}\beta_{j,k}^2\sum_{i=1}^{n}x_{iI_{j,k}}^2$$

$$= n - 1 - \mathbb{1}(k \geq 2)2\sum_{j=2}^{d_k}\beta_{j,k}\sum_{i=1}^{n}y_i x_{iI_{j,k}} + n\beta_{1,k}^2$$

$$+ \mathbb{1}(k \geq 2)(n-1)\sum_{j=2}^{d^k}\beta_{j,k}^2,$$

using again $x_{11} = \ldots = x_{n1} = 1$, $\sum_{i=1}^{n}y_i = 0$ and $\sum_{i=1}^{n}x_{ij}^2 = n - 1$ for all $j \in \{2, \ldots, d\}$. We also have

$$\mathbb{1}(k \geq 2)\left((n-1)\sum_{j=2}^{d_k}\beta_{j,k}^2 - 2\sum_{j=2}^{d_k}\beta_{j,k}\sum_{i=1}^{n}y_i x_{iI_{j,k}}\right)$$

$$= \mathbb{1}(k \geq 2)(n-1)\sum_{j=2}^{d_k}\left(\beta_{j,k}^2 - 2\beta_{j,k}\frac{\sum_{i=1}^{n}y_i x_{iI_{j,k}}}{n-1}\right)$$

$$= \mathbb{1}(k \geq 2)(n-1)\sum_{j=2}^{d_k}\left(\beta_{j,k} - \frac{\sum_{i=1}^{n}x_{iI_{j,k}}y_i}{n-1}\right)^2$$

$$- \mathbb{1}(k \geq 2)(n-1)\sum_{j=2}^{d_k}\left(\frac{\sum_{i=1}^{n}x_{iI_{j,k}}y_i}{n-1}\right)^2.$$

Putting this together leads to:

$$\pi(k, \sigma_k, \boldsymbol{\beta}_k \mid \mathbf{y})$$

$$\propto \pi(k)(2\pi)^{d_k/2} \frac{1}{\sigma_k^{n-d_k+1}} \exp\left\{-\frac{n-1}{2\sigma_k^2}\left(1 - \mathbb{1}(k \geq 2) \sum_{j \in I_k\setminus\{1\}} \left(\frac{\sum_{i=1}^n x_{ij}y_i}{n-1}\right)^2\right)\right\}$$

$$\times \frac{1}{\sigma_k\sqrt{2\pi}} \exp\left\{-\frac{n}{2\sigma_k^2}\beta_{1,k}^2\right\}$$

$$\times \left(\mathbb{1}(k=1) + \mathbb{1}(k \geq 2) \prod_{j=2}^{d_k} \frac{1}{\sigma_k\sqrt{2\pi}} \exp\left\{-\frac{n-1}{2\sigma_k^2}\left(\beta_{j,k} - \frac{\sum_{i=1}^n x_{iI_{j,k}}y_i}{n-1}\right)^2\right\}\right).$$

We multiply and divide by the appropriate terms. The only remaining thing to show is that

$$n-1\left(1 - \mathbb{1}(k \geq 2) \sum_{j \in I_k\setminus\{1\}} \left(\frac{\sum_{i=1}^n x_{ij}y_i}{n-1}\right)^2\right) = \|\mathbf{y} - \widehat{\mathbf{y}}_k\|_2^2.$$

Firstly, $n - 1 = \|\mathbf{y}\|_2^2$. Also,

$$\|\mathbf{y}\|_2^2 = \|\mathbf{y} - \widehat{\mathbf{y}}_k + \widehat{\mathbf{y}}_k\|_2^2$$
$$= \|\mathbf{y} - \widehat{\mathbf{y}}_k\|_2^2 + (\mathbf{y} - \widehat{\mathbf{y}}_k)^T\widehat{\mathbf{y}}_k + \widehat{\mathbf{y}}_k^T(\mathbf{y} - \widehat{\mathbf{y}}_k) + \widehat{\mathbf{y}}_k^T\widehat{\mathbf{y}}_k.$$

We know that $(\mathbf{y} - \widehat{\mathbf{y}}_k)^T\widehat{\mathbf{y}}_k = \widehat{\mathbf{y}}_k^T(\mathbf{y} - \widehat{\mathbf{y}}_k) = 0$ because $\mathbf{y} - \widehat{\mathbf{y}}_k$ is the vector of residuals which is orthogonal to $\widehat{\mathbf{y}}_k$. Finally,

$$\widehat{\mathbf{y}}_k^T\widehat{\mathbf{y}}_k = (\mathbf{X}_k\widehat{\boldsymbol{\beta}}_k)^T\mathbf{X}_k\widehat{\boldsymbol{\beta}}_k = \widehat{\boldsymbol{\beta}}_k^T\mathbf{X}_k^T\mathbf{X}_k\widehat{\boldsymbol{\beta}}_k = (n-1)\|\widehat{\boldsymbol{\beta}}_k\|_2^2$$

$$= (n-1)\mathbb{1}(k \geq 2) \sum_{j \in I_k\setminus\{1\}} \left(\frac{\sum_{i=1}^n x_{ij}y_i}{n-1}\right)^2,$$

where $\mathbf{X}_k$ is the design matrix associated with Model $k$.     □

*Proof of Proposition 2.2.* As explained in Green (1995), it suffices to separately verify that the probability to go from a set $A$ to a set $B$ is equal to the probability to go from $B$ to $A$ when updating the parameters and when switching models, for accepted movements and for any appropriate $A, B$.

When updating the parameters, the probability to go from a set $A$ to a set $B$ is given by

$$\int_A \pi(k, \sigma_k, \boldsymbol{\beta}_k \mid \mathbf{y})g(1) \int_B \prod_{i=1}^{1+d_k} \varphi_i(w_i \mid k, (\sigma_k, \boldsymbol{\beta}_k)_i, \ell_k)$$

$$\times \left(1 \wedge \frac{(1/w_1)f(\mathbf{y} \mid k, \mathbf{w}_k)}{(1/\sigma_k)f(\mathbf{y} \mid k, \sigma_k, \boldsymbol{\beta}_k)}\right) d\mathbf{w}_k \, d(\sigma_k, \boldsymbol{\beta}_k).$$

Using Fubini's theorem, this probability is equal to

$$\int_B \pi(k, \mathbf{w}_k \mid \mathbf{y})g(1) \int_A \prod_{i=1}^{1+d_k} \varphi_i((\sigma_k, \boldsymbol{\beta}_k)_i \mid k, w_i, \ell_k)$$

$$\times \left(1 \wedge \frac{(1/\sigma_k)f(\mathbf{y} \mid k, \sigma_k, \boldsymbol{\beta}_k)}{(1/w_1)f(\mathbf{y} \mid k, \mathbf{w}_k)}\right) d(\sigma_k, \boldsymbol{\beta}_k) \, d\mathbf{w}_k,$$

which is the probability to go from $B$ to $A$. Note that this is valid for all $k \in \{1, \ldots, K_{\max}\}$.

The probability to switch from Model $k \in \{1, \ldots, K_{\max} - 1\}$, where the parameters are in the set $A$, to Model $k + 1$, where the parameters are in the set $A' \times B$ (the set $A'$ is a modified version of $A$ to account for the addition of $\mathbf{c}_{k+1}$), is given by

$$
\int_A \pi(k, \sigma_k, \boldsymbol{\beta}_k \mid \mathbf{y}) g(2) \int_B q_{k+1}(u_{k+1})
$$
$$
\times \left( 1 \wedge \frac{\pi(k+1) f(\mathbf{y} \mid k+1, (\sigma_k, \boldsymbol{\beta}_k) + \mathbf{c}_{k+1}, u_{k+1})}{\pi(k) f(\mathbf{y} \mid k, \sigma_k, \boldsymbol{\beta}_k) q_{k+1}(u_{k+1})} \right) du_{k+1} \, d(\sigma_k, \boldsymbol{\beta}_k).
$$

After the change of variables $(\sigma_{k+1}, \boldsymbol{\beta}_{k+1}) = ((\sigma_k, \boldsymbol{\beta}_k) + \mathbf{c}_{k+1}, u_{k+1})$, we have

$$
\int_{A' \times B} \pi(k, (\sigma_{k+1}, \boldsymbol{\beta}_{k+1}^-) - \mathbf{c}_{k+1} \mid \mathbf{y}) g(2) q_{k+1}(\beta_{d_{k+1}, k+1})
$$
$$
\times \left( 1 \wedge \frac{\pi(k+1) f(\mathbf{y} \mid k+1, \sigma_{k+1}, \boldsymbol{\beta}_{k+1})}{\pi(k) f(\mathbf{y} \mid k, (\sigma_{k+1}, \boldsymbol{\beta}_{k+1}^-) - \mathbf{c}_{k+1}) q_{k+1}(\beta_{d_{k+1}, k+1})} \right) d(\sigma_{k+1}, \boldsymbol{\beta}_{k+1}).
$$

This last probability is equal to

$$
\int_{A' \times B} \pi(k+1, \sigma_{k+1}, \boldsymbol{\beta}_{k+1} \mid \mathbf{y}) g(3)
$$
$$
\times \left( 1 \wedge \frac{\pi(k) f(\mathbf{y} \mid k, (\sigma_{k+1}, \boldsymbol{\beta}_{k+1}^-) - \mathbf{c}_{k+1}) q_{k+1}(\beta_{d_{k+1}, k+1})}{\pi(k+1) f(\mathbf{y} \mid k+1, \sigma_{k+1}, \boldsymbol{\beta}_{k+1})} \right) d(\sigma_{k+1}, \boldsymbol{\beta}_{k+1}),
$$

which is the probability to switch from Model $k + 1$, where the parameters are in the set $A' \times B$, to Model $k$, where the parameters are in the set $A$.

Therefore, the Markov chain $\{(K, \sigma_K, \boldsymbol{\beta}_K)(m) : m \in \mathbb{N}\}$ satisfies the reversibility condition with respect to the posterior. $\qquad \square$

## 7.4   List of the explanatory variables used in Section 6

| Name | Ticker symbol |
|---|---|
| Artis Real Estate Investment Trust | AX-UN.TO |
| Asanko Gold Inc. | AKG.TO |
| Bonterra Energy Corp. | BNE.TO |
| Canadian Imperial Bank Of Commerce | CM.TO |
| CI Financial Corp. | CIX.TO |
| Celestica Inc. Subordinate Voting Shares | CLS.TO |
| DHX Media Ltd. | DHX-B.TO |
| Dominion Diamond Corporation | DDC.TO |
| Gildan Activewear Inc. | GIL.TO |
| Husky Energy Inc. | HSE.TO |
| iPath Bloomberg Sugar Subindex | SGG |
| iShares MSCI Japan | EWJ |
| iShares 20+ Year Treasury Bond | TLT |
| Laurentian Bank of Canada | LB.TO |
| Parkland Fuel Corporation | PKI.TO |
| United States Oil Fund LP | USO |
| Vermilion Energy Inc. | VET.TO |
| Volume of the S&P 500 | N/A |

**Table 2.** Names of the companies, funds, and financial indicators used as explanatory variables in the analysis in Section 6 of our paper, with their ticker symbol (if available)